

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**Method and Apparatus for Retrieving and
Processing Data**

Inventors:

Roy Messing

Jeremy Sokolic

Venkatachari Dilip

Sanjeev Dheer

ATTORNEY'S DOCKET NO. CE1-005US

1 **TECHNICAL FIELD**

2 The present invention relates to the retrieval and processing of data
3 collected from web pages and/or other data sources.
4

5 **BACKGROUND**

6 Individuals, businesses, and other organizations typically maintain one or
7 more financial accounts at one or more financial institutions. Financial institutions
8 include, for example, banks, savings and loans, credit unions, mortgage
9 companies, lending companies, and stock brokers. A customer's financial
10 accounts may include asset accounts (such as savings accounts, checking accounts,
11 certificates of deposit (CDs), mutual funds, bonds, and equities) and debt accounts
12 (such as credit card accounts, mortgage accounts, home equity loans, overdraft
13 protection, and other types of loans).

14 Many financial institutions allow customers to access information regarding
15 their accounts via the Internet or other remote connection mechanism (often
16 referred to as "online banking"). Typically, the customer navigates, using a web
17 browser application, to a web site maintained by the financial institution. The web
18 site allows the customer to login by entering a user identification and an associated
19 password. If the financial institution accepts the user identification and password,
20 the customer is permitted to access information (e.g., account holdings and
21 account balances) regarding the financial accounts maintained at that financial
22 institution.

23 Similarly, other organizations and institutions allow customer access to
24 other types of accounts, such as email accounts, award (or reward) accounts,
25 online bill payment accounts, etc. A user may navigate a web site or other

1 information source to receive status information regarding one or more of their
2 accounts.

3 Certain application programs are able to extract data from web pages based
4 on a previously defined layout of information on the web pages. For example, an
5 account balance may be positioned in a particular location of a specific web page.
6 The application program extracts the account balance data from that particular
7 location to obtain a customer's current account balance. However, if the layout of
8 the web page is modified, the previously defined layout of information on the web
9 page is not accurate and the application program cannot properly extract the
10 desired data from the web page.

11 The systems and methods described herein addresses these and other
12 problems by providing a mechanism for updating the manner in which data is
13 extracted from web pages when a web page layout is changed.

14 15 16 SUMMARY

17 The systems and methods described herein automatically extract data from
18 web pages and other data sources associated with various institutions. The data is
19 extracted from a data source, such as a web page using a data harvesting script or
20 other data extraction/data acquisition routine. The extracted data is stored in a
21 database using a standard format. If the layout of data on a particular web page
22 changes, a copy of the web page is captured and stored for future analysis when
23 updating one or more data extraction procedures (e.g., data harvesting scripts).
24 Personal or confidential information is deleted from the captured web page before
25 storing the captured web page.

1 A particular embodiment captures a web page from an institution's web
2 site. Data is extracted from the web page using a data harvesting script. The
3 extracted data is then normalized and stored in a database.

4 In another embodiment, a web page is captured from a web site. An
5 attempt is made to extract data from the web page using a data harvesting script.
6 If data cannot be extracted from the web page, personal information is removed
7 from the captured web page and the captured web page (without the personal
8 information) is stored.

9 10 11 **BRIEF DESCRIPTION OF THE DRAWINGS**

12 Fig. 1 illustrates an exemplary network environment in which various
13 servers, computing devices, and a financial analysis system exchange data across a
14 network, such as the Internet.

15 Fig. 2 is a block diagram showing exemplary components and modules of a
16 financial analysis system.

17 Figs. 3A and 3B are flow diagrams illustrating procedures for retrieving
18 data from an HTML screen and another data source.

19 Figs. 4 and 5 illustrate exemplary web pages associated with a particular
20 financial institution.

21 Fig. 6 is a flow diagram illustrating a procedure for retrieving financial data
22 and adjusting a data harvesting script if the financial data layout has changed.

23 Fig. 7 is a block diagram showing pertinent components of a computer in
24 accordance with the invention.
25

DETAILED DESCRIPTION

The system and methods described herein are capable of automatically extracting data from web pages or other data sources associated with one or more accounts or institutions, such as financial accounts or financial institutions. A particular web page or data source may contain account information associated with a customer of a particular institution. If an error occurs when attempting to extract data from a web page, a copy of the web page is saved for future analysis in determining the cause of the error and creating a new procedure for extracting data from the web page. When saving a copy of the web page for future analysis, confidential information is removed before storing the web page, thereby reducing the possibility of inadvertently exposing confidential information contained in the web page.

As used herein, the terms “account holder”, “customer”, “user”, and “client” are interchangeable. “Account holder” refers to any person having access to an account. A particular account may have multiple account holders (e.g., a joint checking account having husband and wife as account holders or a corporate account identifying several corporate employees as account holders).

Various financial account and financial institution examples are provided herein for purposes of explanation. However, the methods and procedures described herein can be applied to any type of transaction involving any type of account. For example, a data aggregation system may aggregate data from multiple sources, such as multiple financial accounts, multiple email accounts, multiple online award (or reward) accounts, multiple news headlines, and the like. Similarly, the data retrieval and data processing systems and methods discussed herein may be applied to collect data from any type of account containing any type

1 of data. Thus, the methods and systems described herein can be applied to a data
2 aggregation system or any other account management system instead of the
3 financial analysis system discussed in the examples provided herein.

4 Fig. 1 illustrates an exemplary network environment 100 in which various
5 servers, computing devices, and a financial analysis system exchange data across a
6 network, such as the Internet. The network environment of Fig. 1 includes
7 multiple financial institution servers 102, 104, and 106 coupled to a data
8 communication network 108, such as the Internet. Data communication network
9 108 may be any type of data communication network using any network topology
10 and any communication protocol. Further, network 108 may include one or more
11 sub-networks (not shown) which are interconnected with one another.

12 A client computer 110 and a financial analysis system 112 are also coupled
13 to network 108. Financial analysis system 112 includes a database 114 that stores
14 various data collected and generated by the financial analysis system. Financial
15 analysis system 112 performs various account analysis and data analysis functions,
16 as discussed in greater detail below.

17 Client computer 110 and financial analysis system 112 may be any type of
18 computing device, such as a desktop computer, a laptop computer, a palmtop
19 computer, a personal digital assistant (PDA), a cellular phone, or a set top box.
20 Client computer 110 communicates with one or more financial institution servers
21 102-106 to access, for example, information about the financial institution and
22 various user accounts that have been established at the financial institution. Each
23 of the financial institution servers 102-106 is typically associated with a particular
24 financial institution and store data for that financial institution.
25

1 The communication links shown between network 108 and the various
2 devices (102, 104, 106, 110, and 112) shown in Fig. 1 can use any type of
3 communication medium and any communication protocol. For example, one or
4 more of the communication links shown in Fig. 1 may be a wireless link (e.g., a
5 radio frequency (RF) link or a microwave link) or a wired link accessed via a
6 public telephone system or another communication network.

7 Fig. 2 is a block diagram showing exemplary components and modules of
8 financial analysis system 112. A communication interface 202 allows the
9 financial analysis system 112 to communicate with other devices, such as one or
10 more financial institution servers. In one embodiment, communication interface
11 202 is a network interface to a local area network (LAN), which is coupled to
12 another data communication network, such as the Internet.

13 A database control module 204 allows financial analysis system 112 to
14 store data to database 114 and retrieve data from the database. Financial analysis
15 system 112 also stores various financial institution data 206, which may be used to
16 locate and communicate with various financial institution servers. Financial
17 institution data 206 includes, for example, account balance information,
18 transaction descriptions, transaction amounts, and security holdings.

19 A variety of data harvesting scripts 208 are also maintained by financial
20 analysis system 112. For example, a separate data harvesting script 208 may be
21 maintained for each financial institution from which data is extracted. Data
22 harvesting (also referred to as “screen scraping”) is a process that allows, for
23 example, an automated script to retrieve data from one or more web pages
24 associated with a web site. Data harvesting may also include retrieving data from
25

1 a data source using any data acquisition or data retrieval procedure. Additional
2 details regarding data harvesting and data harvesting scripts are provided below.

3 Financial analysis system 112 includes a data capture module 210 and a
4 data extraction module 214. The data capture module 210 captures data (such as
5 web pages or OFX data) from one or more data sources. The data extraction
6 module 214 retrieves (or extracts) data from the captured web pages or other data
7 sources. The data extraction module 214 may use one or more data harvesting
8 scripts 208 to retrieve data from a web page. A personal information filter
9 module 212 removes confidential information from a web page. Thus, the
10 majority of the content of the web page can be stored for future access without
11 risking exposure of an account holder's confidential information.

12 Data capture module 210 may also retrieve data from sources other than
13 web pages. For example, data capture module 210 can retrieve data from a source
14 that supports the Open Financial Exchange (OFX) specification or the Quicken
15 Interchange Format (QIF). OFX is a specification for the electronic exchange of
16 financial data between financial institutions, businesses and consumers via the
17 Internet. OFX supports a wide range of financial activities including consumer
18 and business banking, consumer and business bill payment, bill presentment, and
19 investment tracking, including stocks, bonds, mutual funds, and 401(k) account
20 details. QIF is a specially formatted text file that allows a user to transfer Quicken
21 transactions from one Quicken account register into another Quicken account
22 register or to transfer Quicken transactions to or from another application that
23 supports the QIF format.

24 A failure analysis module 218 in financial analysis system 212 analyzes the
25 failure of a data harvesting script and determines why the script failed. For

1 example, if a web page is redesigned by a financial institution, a data harvesting
2 script that has not been updated to reflect the new web page design may not
3 operate properly. In this situation, the information sought by the data harvesting
4 script may have been moved to a different location on the new web page. The
5 failure analysis module 218 assists a user in identifying the reason for the script
6 failure. A script editing module 216 assists a user in editing a data harvesting
7 script to function properly with a new web page design.

8 Figs. 3A and 3B are flow diagrams illustrating procedures for retrieving
9 data from an HTML screen and another data source. Specifically, Fig. 3A is a
10 flow diagram illustrating a procedure 300 for retrieving data from an HTML
11 screen. Initially, the procedure 300 captures an HTML (HyperText Markup
12 Language) screen from a financial institution web site (block 302). For example,
13 the HTML screen may be a web page associated with the financial institution.
14 Data is then extracted from the HTML screen using a data harvesting script (block
15 304). The extracted data is then normalized (block 306), which refers to the
16 process of arranging the extracted data into a standard format such that data
17 collected from a variety of different web pages is arranged (or normalized) into the
18 same format. The normalized data is then stored in the database (e.g., database
19 114 in Fig. 1) for future reference (block 308).

20 The normalizing of data is useful when collecting data from multiple
21 sources (e.g., multiple financial institutions). Each financial institution may use
22 different terms for the same type of data. For example, one financial institution
23 may use the term “buy” while another financial institution uses the term
24 “purchase” for the same type of transaction. By normalizing the data, a single
25 database can be used to store financial information related to multiple different

1 financial institutions. Thus, various financial analysis tools and procedures can be
2 used to analyze data across multiple financial institutions or other data sources.

3 As mentioned above, data harvesting (or screen scraping) is a process that
4 allows an automated script to retrieve data from a web site and store the retrieved
5 data in a database. The data harvesting scripts are capable of navigating web sites
6 and capturing individual HTML pages. For example, JavaScript and images may
7 be removed from the HTML pages or converted into HTML text if it contains
8 account information. A parser then converts the HTML data into a field-delimited
9 XML format. The XML data communicates with enterprise java beans (EJBs)
10 through an XML converter. EJBs perform a series of SQL queries that populate
11 the data into the database. The success of a particular data harvesting process is
12 related to the layout of the web site being harvested in two important ways: 1) the
13 data harvesting script must navigate to the correct HTML page, and 2) the parser
14 must know which cells in the HTML tables contain specific data items.

15 Fig. 3B is a flow diagram illustrating a procedure 350 for retrieving and
16 processing data from a data source (other than an HTML screen). The data source
17 may be, for example, a financial institution or other provider of financial data.
18 The data source may also be referred to as a "file download source" or a "data
19 download source". The data source may communicate data using the OFX
20 standard, the QIF format, or some other data format. The procedure 350 begins by
21 retrieving data from a data source (block 352). The procedure identifies data of
22 interest from the retrieved data (block 354). The data of interest may be, for
23 example, data associated with a particular customer's accounts. The identified
24 data is then normalized (block 356) and stored in the database (block 358). The
25

1 database may contain data related to other customers and/or data collected from
2 other sources (such as HTML screens).

3 Figs. 4 and 5 illustrate exemplary web pages 400 and 500, respectively,
4 associated with a particular financial institution. A particular data harvesting
5 script may look for specific text on the web page to confirm that the script has
6 navigated to the correct site. For example, to ensure that the Vanguard screen
7 scraping script has navigated to the “Quick Links” web page, the script looks for
8 the phrase “Common Tasks” in row 1 of table 1 (see the portion of the web page
9 surrounded by a ring 410). If this phrase is found, the script can then navigate to
10 row two and select the “Access my Accounts” link that takes the script to a secure
11 login page (e.g., an HTTPS login page). If the script cannot locate the phrase
12 “Common Tasks” it will generate an exception error and stop running.

13 Once the script has found the correct page, pattern matches are used by the
14 parser to determine the appropriate cell from which to retrieve specific data items.
15 For example, once the data harvesting script has navigated to the “Account
16 Values” page (shown in Fig. 5), the script identifies the correct row from which to
17 retrieve data by pattern matching a combination of the fund/account number and
18 the fund name in columns one and two. The script also matches the column
19 header name and then moves down the column to the appropriate row in the
20 column. In this example, the parser will populate the data field “Account Value”
21 with the data in the cell in row one and column five. This account value
22 information is highlighted by a ring 510 in Fig. 5.

23 Fig. 6 is a flow diagram illustrating a procedure 600 for retrieving financial
24 data and adjusting a data harvesting script if the financial data layout has changed.
25 Initially, the procedure 600 captures a financial institution screen shot (block 602).

1 For example, a screen shot associated with a particular financial institution web
2 page or web site. Next, the procedure removes personal and/or confidential
3 information from the screen shot (block 604). Example personal and/or
4 confidential information that is removed includes customer name, address,
5 telephone number, email address, and social security number.

6 The procedure 600 then identifies and sorts all failed updates (block 606).
7 A failed update may occur when a data harvesting script attempts to update a
8 user's account information but the layout of the financial institution's web pages
9 have changed. The procedure may search the database for all failed updates by
10 error code (error codes are discussed in greater detail below). The results of the
11 search are provided to one or more individuals responsible for updating screen
12 scraping scripts. Next, bugs are reported and assigned to a particular individual or
13 group for processing (block 608).

14 At block 610, a user accesses the HTML data (i.e., the screen shot captured
15 from the financial institution) to repair the scripts that are not functioning properly.
16 The procedure then continues to block 612, which captures the next financial
17 institution screen shot. The procedure returns to block 604 to remove personal
18 information from the captured screen shot.

19 When a data harvesting script is unable to access a particular web page (or
20 web site) or is unable to locate information on the web page an error occurs. The
21 data harvesting script contains error detection mechanisms that identify errors and
22 generate one or more error codes associated with the identified errors. Each error
23 has an associated error code that identifies the particular error. Table 1 below
24 identifies several example error codes as well as a corresponding title and
25 description of the error that occurred.

TABLE 1

Error Code	Title	Description
100	Web Page Modified	Unable to retrieve account information from financial institution web page due to changes in web page.
101	Time Out	Unable to retrieve account information due to high network traffic.
102	Connection Failed	Unable to retrieve account information due to network connection problems.
103	Web Site Unavailable	Unable to retrieve account information because the financial institution web site is not available.
104	Login Failure	Unable to retrieve account information because the username/password combination provided by user failed.

Different actions may be performed depending on the error detected. For example, if the web page has been modified, the screen shot of the modified web page is provided to one or more individuals to analyze and update the corresponding data harvesting script to properly extract data from the modified web page. If the error indicates a failed network connection, the financial analysis system may attempt to retrieve the desired web pages at a later time. If the error indicates that the username and/or password provided by the user is incorrect, the financial analysis system may request the user verify the username and password associated with the account being accessed.

The error codes may be processed by an automated error handling routine to notify the proper individual, or group of individuals, of the error. For example, a database error may be automatically routed to a group of individuals responsible for managing the database. Other error codes may indicate a problem with the

1 information provided by the user. These error codes, such as an invalid password
2 to access a user account, result in sending an error notice to the user, but do not
3 represent a problem with the financial analysis system.

4 Fig. 7 is a block diagram showing pertinent components of a computer 700
5 in accordance with the invention. A computer such as that shown in Fig. 7 can be
6 used, for example, to perform various procedures such as those discussed herein.
7 Computer 700 can also be used to access a web site or other computing facility to
8 access various financial information. The computer shown in Fig. 7 can function
9 as a server, a client computer, or a financial analysis system, of the types discussed
10 herein.

11 Computer 700 includes at least one processor 702 coupled to a bus 704 that
12 couples together various system components. Bus 704 represents one or more of
13 any of several types of bus structures, such as a memory bus or memory controller,
14 a peripheral bus, and a processor or local bus using any of a variety of bus
15 architectures. A random access memory (RAM) 706 and a read only memory
16 (ROM) 708 are coupled to bus 704. Additionally, a network interface 710 and a
17 removable storage device 712, such as a floppy disk or a CD-ROM, are coupled to
18 bus 704. Network interface 710 provides an interface to a data communication
19 network such as a local area network (LAN) or a wide area network (WAN) for
20 exchanging data with other computers and devices. A disk storage 714, such as a
21 hard disk, is coupled to bus 704 and provides for the non-volatile storage of data
22 (e.g., computer-readable instructions, data structures, program modules and other
23 data used by computer 700). Although computer 700 illustrates a removable
24 storage 712 and a disk storage 714, it will be appreciated that other types of
25 computer-readable media which can store data that is accessible by a computer,

1 such as magnetic cassettes, flash memory cards, digital video disks, and the like,
2 may also be used in the exemplary computer.

3 Various peripheral interfaces 716 are coupled to bus 704 and provide an
4 interface between the computer 700 and the individual peripheral devices.
5 Exemplary peripheral devices include a display device 718, a keyboard 720, a
6 mouse 722, a modem 724, and a printer 726. Modem 724 can be used to access
7 other computer systems and devices directly or by connecting to a data
8 communication network such as the Internet.

9 A variety of program modules can be stored on the disk storage 714,
10 removable storage 712, RAM 706, or ROM 708, including an operating system,
11 one or more application programs, and other program modules and program data.
12 A user can enter commands and other information into computer 700 using the
13 keyboard 720, mouse 722, or other input devices (not shown). Other input devices
14 may include a microphone, joystick, game pad, scanner, satellite dish, or the like.

15 Computer 700 may operate in a network environment using logical
16 connections to other remote computers. The remote computers may be personal
17 computers, servers, routers, or peer devices. In a networked environment, some or
18 all of the program modules executed by computer 700 may be retrieved from
19 another computing device coupled to the network.

20 Typically, the computer 700 is programmed using instructions stored at
21 different times in the various computer-readable media of the computer. Programs
22 and operating systems are often distributed, for example, on floppy disks or CD-
23 ROMs. The programs are installed from the distribution media into a storage
24 device within the computer 700. When a program is executed, the program is at
25 least partially loaded into the computer's primary electronic memory. As

1 described herein, the invention includes these and other types of computer-
2 readable media when the media contains instructions or programs for
3 implementing the steps described below in conjunction with a processor. The
4 invention also includes the computer itself when programmed according to the
5 procedures and techniques described herein.

6 For purposes of illustration, programs and other executable program
7 components are illustrated herein as discrete blocks, although it is understood that
8 such programs and components reside at various times in different storage
9 components of the computer, and are executed by the computer's processor.
10 Alternatively, the systems and procedures described herein can be implemented in
11 hardware or a combination of hardware, software, and/or firmware. For example,
12 one or more application specific integrated circuits (ASICs) can be programmed to
13 carry out the systems and procedures described herein.

14 Although the description above uses language that is specific to structural
15 features and/or methodological acts, it is to be understood that the invention
16 defined in the appended claims is not limited to the specific features or acts
17 described. Rather, the specific features and acts are disclosed as exemplary forms
18 of implementing the invention.